

Durham, North Carolina, U.S.A.

GPO PRICE	\$
CFSTI PRICE(S)	\$

Hard copy (HC) 3.00

Microfiche (MF) 165

Kendall & Stuart (1961, Vol. II, pp. 87 et seq.) cite the work of Lloyd (1952) and Downton (1953) in estimating the scale and location of a random variable by forming the best linear unbiassed estimate based on its order statistics. The purpose of this note is to point out the rather curious consequences which ensue when this sort of technique is applied to the estimation of the scale of an exponential distribution of known location, a situation which arises not infrequently in practice. The result which emerges is that the information relevant to the scale is concentrated markedly in the upper portion of the sample, there being a logarithmic type of singularity such that the upper half of the sample contains 98.0% of the total information in the sample; the upper 10%, 68.9% of the information; and the upper 1%, 22.4%.

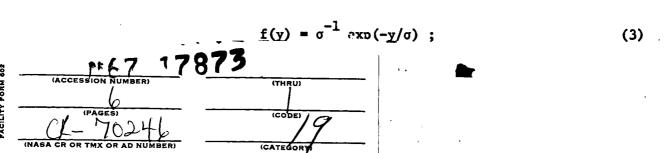
We assume that we have an ordered sample

$$\underline{y}_{(1)} \leq \underline{y}_{(2)} \leq \cdots \leq \underline{y}_{(n)} \tag{1}$$

of an exponential variate

$$\underline{y} \sim \exp 1(\sigma)$$
, (2)

by which we mean that the frequency function of y is



and we write

$$\underline{y} \equiv \sigma \underline{z}$$
, $\underline{z} \wedge \exp 1(1)$ (4)

and

$$\underline{y}_{(r)} = \sigma \underline{z}_{(r)} \quad (\underline{r} = 1, 2, ..., \underline{n}) . \tag{5}$$

Then the column vector

$$z = \{\underline{z}_{(n)}, \underline{z}_{(n-1)}, \dots, \underline{z}_{(n-k+1)}\}$$
 (6)

(we assume henceforth that \underline{k} is a fixed integer between 1 and n) has, as Kendall & Stuart (op. $\underline{\text{cit.}}$, pp. 96, 97, problem 19.11) remark, a mean given by

$$g = \{ \sum_{j=1}^{n} \frac{1}{j}, \sum_{j=1}^{n} \frac{1}{j}, \dots, \sum_{k=1}^{n} \frac{1}{j} \}$$
 (7)

and a variance equal to

These results are easy consequences of the fact that the ordered sample (5) of the z's may also be generated by the equivalent sampling scheme of choosing

$$\underline{\mathbf{u}}_{1} = \underline{\mathbf{z}}_{(1)} \sim \exp(1/\underline{\mathbf{n}}) \tag{9}$$

and

$$\frac{\mathbf{u}}{\mathbf{j}} = \underline{z}_{(j)} - \underline{z}_{(j-1)} \sim \exp[1/(\underline{n}-\underline{j}+1)] \quad (\underline{j} = 2, 3, ..., \underline{n}) , \quad (9')$$

where the $\underline{\mathbf{u}}$'s are all mutually independent; for, this makes the joint frequency function

$$f(\underline{z}_{(1)}, \dots, \underline{z}_{(n)})$$
= $\underline{n} \exp[-\underline{n}\underline{z}_{(1)}] \prod_{j=1}^{n-1} (\underline{n}-\underline{j})\exp\{-(\underline{n}-\underline{j})[\underline{z}_{(j+1)} - \underline{z}_{(j)}]\}$
= $\underline{n}! \prod_{j=1}^{n} \exp[-\underline{z}_{(j)}]$, (10)

which is clearly what it should be.

Corresponding to (6), we write the top \underline{k} members of the ordered sample of the original variate \underline{y} as the column vector

$$y = \{y_{(n)}, y_{(n-1)}, \dots, y_{(n-k+1)}\}$$
 (11)

and standard results from linear estimation then state that for the least squares estimate we have

$$\hat{\sigma} = \underline{S}^{-1} \underline{\alpha} , \underline{V}^{-1} \underline{V} \quad \text{and} \quad \text{var } \hat{\sigma} = \underline{\sigma}^{2} \underline{S} , \qquad (12)$$

where

$$\underline{\mathbf{s}} = \mathbf{g}^{\mathsf{T}} \mathbf{y}^{-1} \mathbf{g} , \qquad (13)$$

a scalar, since we have only a single parameter here.

Now, one may readily check that

$$\chi^{-1} =
\begin{pmatrix}
1^2 & -1^2 & 0 & \cdots & 0 & 0 \\
-1^2 & 1^2 + 2^2 & -2^2 & \cdots & 0 & 0 \\
0 & -2^2 & 2^2 + 3^2 & \cdots & 0 & 0 \\
\vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\
0 & 0 & 0 & \cdots & (\underline{k}-2)^2 + (\underline{k}-1)^2 & -(\underline{k}-1)^2 \\
0 & 0 & 0 & \cdots & -(\underline{k}-1)^2 & (\underline{k}-1)^2 + \frac{1}{p} \frac{1}{\frac{1}{p} \frac{1}{4}^2}
\end{pmatrix}$$
(14)

and from this it follows that

$$\chi'\chi^{-1} = (1, 1, ..., 1, -k+1 + \frac{\sum_{k=1}^{n} \frac{1}{j}}{\sum_{k=1}^{n} \frac{1}{2}})$$
 (15)

and

$$\underline{S} = \underline{k} - 1 + \frac{\left(\sum_{k=1}^{n} \frac{1}{j}\right)^{2}}{\sum_{k=1}^{n} \frac{1}{2}}.$$
 (16)

Thus, the best linear unbiassed estimate of the scale factor, based on the top k members of the sample, is given by

$$\hat{\sigma} = \frac{1}{S} \left\{ \sum_{n-k+1}^{n} y_{(i)} + \left[\frac{\sum_{k=1}^{n} \frac{1}{j}}{\sum_{k=1}^{n} \frac{1}{2}} - k \right] y_{(n-k+1)} \right\}, \qquad (17)$$

a rather curious expression, since, for moderate values of k, the second principal term within the curly brackets completely swamps the first.

When $\underline{k} = \underline{n}$, (16) yields $\underline{S} = \underline{n}$ and (17) reduces to the arithmetic mean, which is, of course, the maximum likelihood estimator based on the entire sample. Therefore, by the second half of (12), the relative efficiency of (17) is $\underline{S/n}$; but, writing $\underline{k} = \theta \underline{n}$ (0 < θ < 1), we have

$$\sum_{k=1}^{n} \frac{1}{j} \sim \int_{\theta n}^{n} \frac{dx}{x} = -\log \theta \quad (n \to \infty)$$
 (18)

and

$$\sum_{k=1}^{n} \frac{1}{2} \sim \int_{\theta n}^{n} \frac{dx}{x^{2}} = \left(\frac{1}{\theta} - 1\right) \frac{1}{n} \quad (n \to \infty) \quad . \tag{19}$$

Hence, (16) yields

$$\underline{S/n} \sim \left\{1 + \frac{(\log \theta)^2}{1 - \theta}\right\} \theta \quad (k = \theta n, n + \infty), \qquad (20)$$

the result alluded to at the outset.

An obvious application of this result is to the estimation of frequency functions which are linear combinations of distinct exponentials:

$$\underline{\mathbf{f}}(\mathbf{y}) = \sum_{i=1}^{\nu} \underline{\mathbf{A}}_{i} \exp(-\underline{\mathbf{y}}/\sigma_{i}) \quad (\sigma_{1} > \sigma_{2} > \dots > \sigma_{\nu}) . \tag{21}$$

Since

$$\underline{\mathbf{f}}(\mathbf{y}) \sim \underline{\mathbf{A}}_1 \exp(-\mathbf{y}/\sigma_1) \quad (\mathbf{y} \rightarrow \infty) ,$$
 (22)

in a sufficiently large sample size the upper extreme values will effectively be governed by (22), and we may use (17) to estimate the dominant "eigenvalue," σ_1 , provided \underline{k} is not taken to be large enough to introduce contributions from the smaller σ 's in (21), a delicate matter.

When it is required to estimate location as well as scale, the behaviour we have been describing no longer obtains.

ACKNOWLEDGMENT

This work was done under grant number NGR 34-001-005 from the National Aeronautics and Space Administration of the U.S. Government.

REFERENCES

- Downton, F. (1953). A note on ordered least-squares estimation. Biometrika, 40, 457.
- Kendall, M. G. and Stuart, A. (1961). The Advanced Theory of Statistics.

 Griffin, London.
- Lloyd, E. H. (1952). Least-squares estimation of location and scale parameters using order statistics. Biometrika, 39, 88.